

SATISFIED USER RATIO PREDICTION WITH SUPPORT VECTOR REGRESSION FOR COMPRESSED STEREO IMAGES

Chunling Fan*, Yun Zhang*, Raouf Hamzaoui[†], Djemel Ziou[‡], Qingshan Jiang*

*Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, China

[†]School of Engineering and Sustainable Development, De Montfort University, UK

[‡]Département Faculté d'informatique, Université de Sherbrooke, Québec, Canada

ABSTRACT

We propose the first method to predict the Satisfied User Ratio (SUR) for compressed stereo images. The method consists of two main steps. First, considering binocular vision properties, we extract three types of features from stereo images: image quality features, monocular visual features, and binocular visual features. Then, we train a Support Vector Regression (SVR) model to learn a mapping function from the feature space to the SUR values. Experimental results on the SIAT-JSSI dataset show excellent prediction accuracy, with a mean absolute SUR error of only 0.08 for H.265 intra coding and only 0.13 for JPEG2000 compression.

Index Terms— Satisfied user ratio, picture-level just noticeable difference, stereo images

1. INTRODUCTION

The Picture-level Just Noticeable Difference (PJND) for a given image and compression scheme is the smallest distortion level that can be perceived by a subject when the image is compared to its compressed versions. If the distortion level resulting from compression is lower than a subject's PJND, then the subject is satisfied with the reconstructed image quality. Due to natural variations in human physiological structure, visual acuity, and visual attention mechanisms, the PJND may differ from one person to another. The Satisfied User Ratio (SUR) for a given distortion level is the fraction of subjects whose PJND is greater than this distortion level. The SUR can be exploited in streaming applications to save bit rate while keeping a good visual quality. The most straightforward way to determine the PJND and SUR for a population is through subjective image quality assessment, where a group of subjects are invited to view images or videos and compare them to their compressed versions. As it is expensive and time consuming to build a large-scale dataset, it is meaningful to build objective models to predict the PJND and SUR for images and videos.

In recent years, the PJND has become a hot topic in image quality assessment. Researchers have proposed a number of works to study the PJND characteristics of the Human

Visual System (HVS). Jin et al. [1] studied the PJND characteristics of 2D images compressed with JPEG and found that there are only four to seven distinct perceptual quality levels for most images. Wang et al. [2] explored the PJND characteristics of 2D videos compressed with H.264/AVC. Wang et al. [3] obtained the PJND and SUR for each video encoded with H.264/AVC by statistical analysis of the collected PJND samples. Huang et al. [4] modeled the PJND of HEVC compressed 2D video as a normal distribution and proposed a Support Vector Regression (SVR)-based model to predict the mean of the distribution. Hadizadeh et al. [5] developed a sparse coding based model to predict whether a given JND-noise-contaminated image is perceptually distinguishable from a reference image. First, they cropped non-overlapping patches from the reference and distorted images and selected patches with lower quality than the average level. Then, these patches were sparsely coded using a learned dictionary and a feature vector was fed to a binary classifier composed of a multi-layer feedforward neural network. Liu et al. [6] proposed a deep learning-based model to predict the PJND of compressed images. They first modeled the prediction task as a multi-classification problem and transformed it into a binary classification problem. As a binary classifier, a deep learning predictor was used. Its task is to predict whether the distorted image is perceptually lossy compared to its reference. In addition to the works that predict the PJND, other works were proposed to predict the SUR. Wang et al. [7] used an SVR to predict the SUR of H.264/AVC compressed video. Fan et al. [8] developed a deep learning-based method to predict the SUR for JPEG compressed 2D images.

Because the real world is three-dimensional, stereo images provide a more realistic visual experience than 2D images. However, models designed for 2D images cannot be directly applied to stereo images because of binocular vision. It is desirable to develop models to predict the SUR for stereo images. In [9], we studied the PJND characteristics of symmetrically and asymmetrically compressed stereo images and generated two PJND-based stereo image datasets. In this paper, we build the first objective model that can predict the SUR for compressed stereo images.

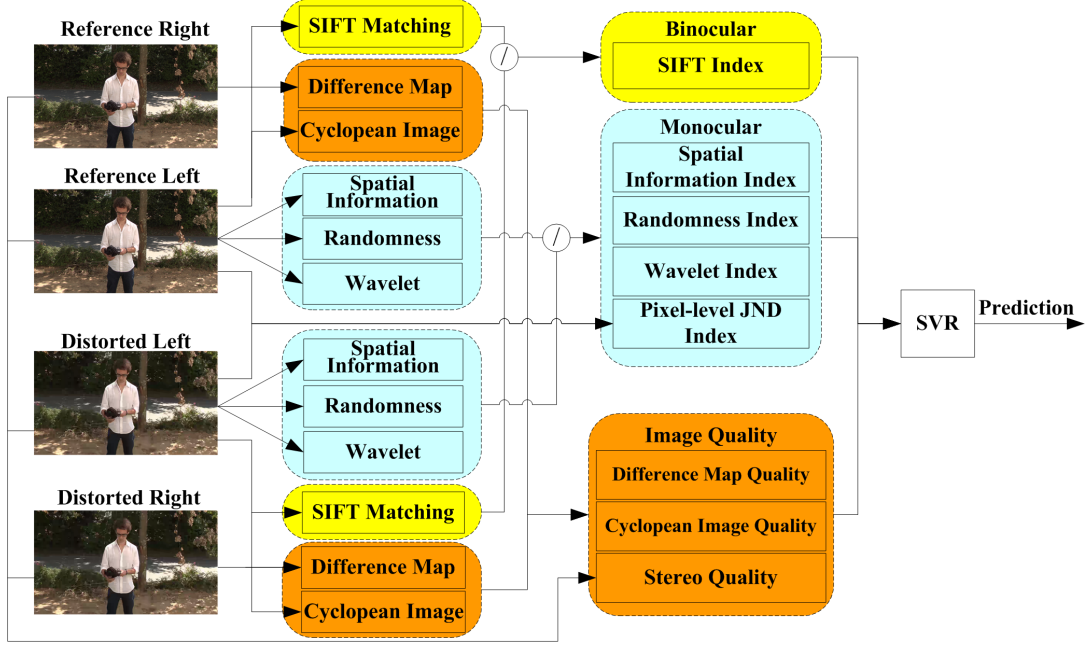


Fig. 1: Proposed architecture for SUR prediction.

2. PROPOSED METHOD

2.1. Problem Modeling

In this paper, we modeled SUR prediction for stereo images as a regression problem. We used an SVR-based model because it is very effective to solve regression problems for high-dimensional features. Let $I_k^L[0]$ and $I_k^R[0]$, $k = 1, \dots, K$ be the left views and right views, respectively, of K pristine stereo images $(I_k^L[0], I_k^R[0])$. Let $(I_k^L[n], I_k^R[n])$, $n = 1, \dots, N$, be the N distorted versions of the pristine stereo image $(I_k^L[0], I_k^R[0])$, where n is the distortion level. Given a pristine stereo image $(I_k^L[0], I_k^R[0])$ and its distorted version $(I_k^L[n], I_k^R[n])$ at distortion level n , we aim to learn an SVR-based model S_θ to predict the SUR value $\text{SUR}_k^n = \text{Prob}[\text{PJND} > n]$. That is

$$S_\theta(f_1, f_2, \dots, f_M) \approx \text{SUR}_k^n, \quad (1)$$

where f_1, f_2, \dots, f_M are the features extracted from the reference and distorted stereo images. The issue is to define the features allowing the prediction of SUR_k^n . If we assume that the PJND samples for each reference image are normally distributed with mean μ and variance σ^2 , then the SUR is given [9] by the Complementary Cumulative Distribution Function (CCDF)

$$\bar{\Phi}(x|\mu, \sigma^2) = 1 - \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(s-\mu)^2}{2\sigma^2}} ds, \quad (2)$$

where μ and σ^2 are to be determined from the PJND samples. In particular, the ground truth SUR values SUR_k^n are given by $\bar{\Phi}(n, |\mu, \sigma^2)$.

Fig. 1 shows our SVR-based model to predict the SUR for compressed stereo images. Our model consists of two main steps. First, we extract features that are sensitive to user perceptual satisfaction from the reference and distorted stereo images. Then we concatenate the extracted features and feed them into an SVR-based model to learn the mapping function from the feature space to SUR values.

2.2. Feature Extraction

Considering the monocular and binocular vision, three kinds of features are chosen: image quality features, monocular visual features, and binocular visual features.

Image Quality Features In image coding, the quality of the compressed image depends on the encoding parameters. For example, in H.265, one encoding parameter is the Quantization Parameter (QP). Specifically, QP=0 gives the highest reconstruction quality, and QP=51 gives the lowest reconstruction quality. We assume that there is a high correlation between the compressed stereo image quality and the SUR. This assumption is realistic because it has been found [9] that there is a high correlation between the SUR and QP for compressed stereo images. We estimated the quality of the distorted stereo images using the Frequency Integrated Peak Signal-to-Noise Ratio (FIPSNR) [10] and denoted it by q_{stereo} . In addition, considering the binocular rivalry and binocular fusion of the HVS, we also estimated the quality of the difference map and the cyclopean image [11] using the Peak Signal-to-Noise Ratio (PSNR) and denoted them by q_{diff} and q_{cyc} , respectively. Finally, we concatenated the three image quality features in a vector $\vec{f}_Q = (q_{\text{stereo}}, q_{\text{diff}}, q_{\text{cyc}})$.

Monocular Visual Features The distortions in an image may not be perceived by the HVS because of the masking effect [12]. Researchers [13, 14] have proposed pixel-level JND models by exploiting the JND characteristics of the HVS. These models mainly focus on luminance and contrast masking effects. First, we used the algorithm in [13] to estimate the pixel-level JND threshold map T for the left view of the reference image. Then, we obtained the ratio of the pixels that are outside the range defined by the threshold T as

$$M_{\text{JND}} = \frac{1}{W \times H} \sum_{i=1}^W \sum_{j=1}^H \Phi(i, j), \quad (3)$$

where

$$\Phi(i, j) = \begin{cases} 1, & \text{if } |R(i, j) - D(i, j)| > T(i, j) \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

W and H are the width and height of the image, R and D are the left view of the reference and distorted stereo image, respectively. We found that when the distortion level increases, the ratio of the pixels that are outside the range defined by the JND threshold increases too. Therefore the ratio M_{JND} was chosen as a feature.

Inspired by the work [15], we extracted spatial randomness to measure the spatial masking effect for stereo images. We first calculated the randomness map [16] for the left view of the reference and distorted images. Then, we extracted a 10-dimensional feature vector from the histogram of the randomness map. To eliminate the effect of the image content, the features $(d_1, d_2, \dots, d_{10})$ from the distorted version were divided by the features $(r_1, r_2, \dots, r_{10})$ from the reference image, giving a randomness feature vector $\vec{M}_{\text{Rand}} = (d_1/r_1, \dots, d_{10}/r_{10})$.

In [9], it was found that there is a high correlation between the Spatial Information (SI) and SUR for stereo images. We calculated SI from the left view of the reference and distorted images, respectively. To eliminate the effect of the image content, we used the ratio of the SI from the distorted image to the SI from the reference image as a feature and denoted it by M_{SI} .

Features extracted from the wavelet transform are often used in image quality assessment. We extracted features from the coefficients obtained after applying the wavelet transform to the left view of the stereo image. Specifically, we extracted the mean and standard deviation of the horizontal, vertical, and diagonal detail coefficients of a level 1 Haar wavelet decomposition, respectively. We concatenated them in a feature vector and denoted it by \vec{M}_{WT} . Finally, we concatenated the four monocular visual features into a single vector $\vec{f}_M = (M_{\text{JND}}, \vec{M}_{\text{Rand}}, M_{\text{SI}}, \vec{M}_{\text{WT}})$.

Binocular Visual Features Binocular vision is the most important characteristic for stereo images compared with 2D images. Because of the horizontal parallax, the two eyes view

an object from slightly different directions, so there exists a positional difference between the two retinal projections. This is known as binocular disparity. The human brain has an ability to deduce depth perception from the binocular disparity. With increasing distortion level, more texture details may be lost which will decrease the information about depth perception. Thus, the number of matching feature points in the left and right views may decrease. The Scale-Invariant Feature Transform (SIFT) [17] descriptor is widely used in stereo matching and depth information extraction [18]. We applied the SIFT to stereo images and used the number of matching points between the left and right view as a binocular visual feature. We found that the number of matching points is related to the image content. To eliminate the effect of the image content, the feature from the distorted version was divided by the feature from the reference image and denoted by f_B .

We concatenated the image quality feature vector, monocular visual feature vector, and binocular visual feature in a single feature vector $\vec{F} = (f_Q, \vec{f}_M, f_B)$. Fig. 1 shows which images were used to compute each feature.

2.3. SUR Prediction Model Learning

Based on the analysis in Section 2.2, we know that the features are related to user satisfaction for compressed stereo images. We aim to learn a mapping function between the feature space and SUR values. An ϵ -SVR-based model [19] was chosen for it is effective for regression of high dimensional features. We consider a training set consisting of K stereo images $(I_k^L[0], I_k^R[0])$, $k = 1, \dots, K$ and their distorted versions $(I_k^L[n], I_k^R[n])$, $k = 1, \dots, K, n = 1, \dots, N$, where n is the number of distortion levels. We assume that for compressed stereo images the PJND samples are normally distributed and use the training data to derive the mean and standard deviation of the distribution. Then, we derive the ground truth SUR as CCDF of this distribution as given in Eq.(2). Next, using the ground truth SUR values, we use the ϵ -SVR model to learn the mapping function between the feature space and the SUR values. Given a test stereo image $(I_k^L[0], I_k^R[0])$ and its N distorted versions $(I_k^L[n], I_k^R[n])$, $n = 1, \dots, N$, we use the learned ϵ -SVR model to predict the N SUR values $\text{SUR}_k^n = \text{Prob}[\text{PJND} > n]$, $n = 1, \dots, N$ given by Eq.(1).

3. EXPERIMENTAL RESULTS AND ANALYSIS

To evaluate the performance of the proposed method, we conducted experiments on the SIAT-JSSI dataset [9]. SIAT-JSSI contains 10 reference stereo images and distorted versions obtained with H.265 intra coding and JPEG2000 compression. For H.265 intra coding, the distorted versions were obtained by varying QP from 1 to 51. For JPEG2000 compression, the distorted versions were obtained by varying the Compression Ratio (CR) from 1 to 300. In total, there are 3510

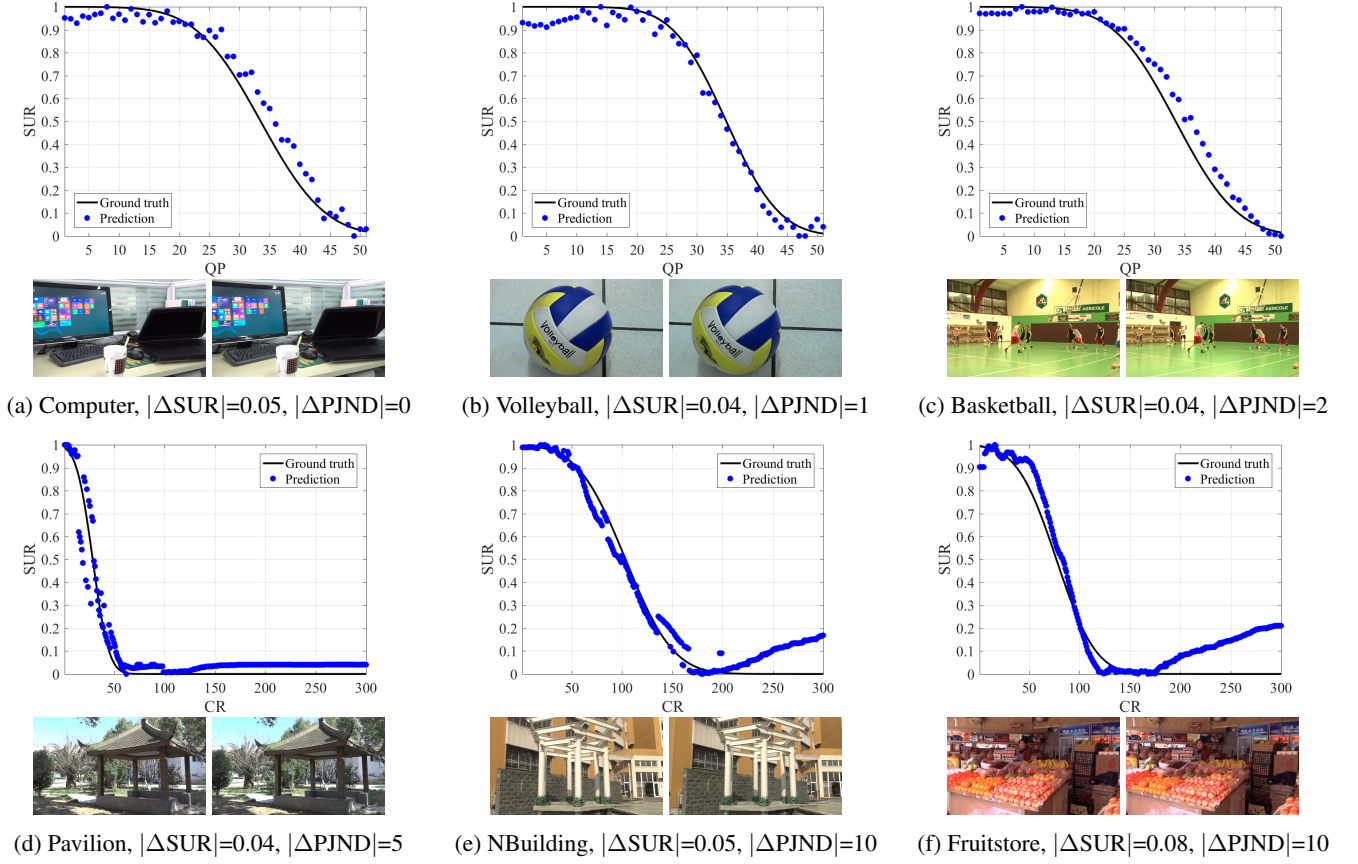


Fig. 2: The three best predicted results. (a)-(c) are for H.265 intra coding, (d)-(f) are for JPEG2000 Compression.

Table 1: Prediction for H.265 intra coding in SIAT-JSSI.

Image	Ground truth		Prediction		$ \Delta PJND $	$ \Delta FIPSNR $	$ \Delta SUR $
	QP	FIPSNR	\widehat{PJND}	\widehat{FIPSNR}			
People	28	54.70	25	56.91	3	2.20	0.07
Basketball	28	57.12	30	55.29	2	1.83	0.04
Newsreport	33	54.35	28	58.28	5	3.92	0.09
Treebranches	27	54.90	24	57.26	3	2.35	0.07
Flower	33	51.84	28	56.16	5	4.32	0.08
Computer	28	56.89	28	56.89	0	0.00	0.05
Volleyball	30	55.51	29	56.35	1	0.84	0.04
Pavilion	26	56.55	34	49.88	8	6.67	0.17
NBuilding	30	55.27	27	57.89	3	2.61	0.06
Fruitstore	25	59.34	28	56.80	3	2.54	0.12
Overall	-	-	-	-	3.30	2.73	0.08

Table 2: Prediction for JPEG2000 compression in SIAT-JSSI.

Image	Ground truth		Prediction		$ \Delta\text{PJND} $	$ \Delta\text{FIPSNR} $	$ \Delta\text{SUR} $
	CR	FIPSNR	$\widehat{\text{PJND}}$	$\widehat{\text{FIPSNR}}$			
People	54	50.85	20	56.76	34	5.91	0.20
Basketball	114	51.41	19	63.90	95	12.49	0.14
Newsreport	156	56.15	84	59.68	72	3.53	0.13
Treebranches	27	52.81	12	57.75	15	4.94	0.08
Flower	45	52.47	15	60.88	30	8.41	0.17
Computer	123	53.17	63	58.11	60	4.94	0.10
Volleyball	170	54.03	62	61.39	108	7.36	0.35
Pavilion	20	56.20	25	54.98	5	1.22	0.04
NBuilding	78	55.32	68	56.65	10	1.34	0.05
Fruitstore	56	56.53	66	54.67	10	1.86	0.08
Overall	-	-	-	-	43.90	5.20	0.13

symmetrically compressed stereo images in the dataset. K -fold ($K = 10$) cross validation was used in our experiment. Specifically, the dataset was split into 10 subsets, where each subset contained one reference image and its distorted versions. Each time, nine subsets were used for training and validating, and the remaining one was used for testing. The reported results are the average of the ten tests. We evaluated our method using three metrics: 1) the absolute PJND error ($|\Delta\text{PJND}|$) between the predicted PJND values and the ground truth PJND values for $\text{SUR} = 0.75$. We obtained the ground truth PJND values from the Gaussian CCDF curves at 0.75 and the predicted PJND values from the predicted SUR values as $\arg \min_{n=1, \dots, N} |\text{SUR}_k^n - 0.75|$, 2) the absolute FIPSNR error ($|\Delta\text{FIPSNR}|$) between the predicted FIPSNR and the ground truth FIPSNR at the PJND corresponding to $\text{SUR} = 0.75$, and 3) the average absolute SUR error ($|\Delta\text{SUR}|$) between the predicted SUR values and the ground truth SUR values over the N distortion levels.

3.1. SUR Prediction for H.265 Intra Coding

In this section, the SUR prediction for H.265 intra coding in SIAT-JSSI is analyzed. Table 1 shows the PJND (as QP), the stereo quality of the distorted image compressed with the PJND measured with FIPSNR, $|\Delta\text{PJND}|$, $|\Delta\text{SUR}|$, and $|\Delta\text{FIPSNR}|$ for the 10 source images. We find that when measured with $|\Delta\text{PJND}|$, the prediction error was zero for one (10%) image and smaller than or equal to 5 for nine (90%) images. When measured with $|\Delta\text{FIPSNR}|$, the prediction error was 0 dB for one (10%) image and smaller than 4.0 dB for eight (80%) images. The mean prediction error for $|\Delta\text{PJND}|$, $|\Delta\text{FIPSNR}|$, $|\Delta\text{SUR}|$ was 3.30, 2.73, and 0.08, respectively. Fig 2 (a)-(c) show the three best predicted results according to $|\Delta\text{SUR}|$. The SUR plots are shown at the top of the figure and the left and right views of the stereo images are shown at the bottom.

3.2. SUR Prediction for JPEG2000 compression

In this section, the SUR prediction for JPEG2000 compression in SIAT-JSSI is analyzed. Table 2 shows the PJND (as CR), the stereo quality of the distorted image compressed with the PJND measured with FIPSNR, $|\Delta\text{PJND}|$, $|\Delta\text{SUR}|$, and $|\Delta\text{FIPSNR}|$ for the 10 source images. We find that when measured with $|\Delta\text{PJND}|$, the prediction error was less than or equal to 10 for three (30%) images. When measured with $|\Delta\text{SUR}|$, the prediction error was smaller than or equal to 0.1 for five (50%) images. The mean prediction error for $|\Delta\text{PJND}|$, $|\Delta\text{FIPSNR}|$, and $|\Delta\text{SUR}|$ was 43.90, 5.20, and 0.13, respectively. Fig 2 (d)-(f) show the three best predicted results according to $|\Delta\text{PJND}|$.

4. CONCLUSIONS

We proposed the first method to predict the SUR for compressed stereo images. Both the monocular and binocular vision properties were considered in our method. Image quality features, monocular visual features, and binocular visual features were extracted and concatenated. Predicting the SUR curve for a stereo image allows us to determine the maximum distortion level that cannot be perceived by any given percentage of the population. This result can be exploited in 3D streaming applications to save bit rate while guaranteeing high visual quality.

Acknowledgments

This work was supported in part by the NSFC under Grant 61871372, Guangdong NSF for Distinguished Young Scholar under Grant 2016A030306022, Guangdong Provincial Science and Technology Development under Grant 2017B010110014, Guangdong R&D Grants 2019B010137002 and 2018A030313943, Shenzhen Research Foundation Grant JCYJ20180302145645821, Shenzhen In-

ternational Collaborative Research Project under Grant GJHZ20170314155404913, Shenzhen Science and Technology Program under Grant JCYJ20170811160212033, Guangdong International Science and Technology Cooperative Research Project under Grant 2018A050506063, Membership of Youth Innovation Promotion Association, CAS under Grant 2018392.

5. REFERENCES

- [1] Lina Jin, Joe Yuchieh Lin, Sudeng Hu, Haiqiang Wang, Ping Wang, Ioannis Katsavounidis, Anne Aaron, and C. C. Jay Kuo, "Statistical study on perceived jpeg image quality via MCL-JCI dataset construction and analysis," *Electronic Imaging*, vol. 2016, no. 13, pp. 1–9, 2016.
- [2] Haiqiang Wang, Weihao Gan, Sudeng Hu, Joe Yuchieh Lin, Lina Jin, Longguang Song, Ping Wang, Ioannis Katsavounidis, Anne Aaron, and C C Jay Kuo, "MCL-JCV: a JND-based H.264/AVC video quality assessment dataset," *the 2016 IEEE International Conference on Image Processing*, pp. 1509–1513, 2016.
- [3] Haiqiang Wang, Ioannis Katsavounidis, Jiantong Zhou, Jeonghoon Park, Shawmin Lei, Xin Zhou, Manon Pun, Xin Jin, Ronggang Wang, Xu Wang, Yun Zhang, Jiwu Huang, Sam Kwong, and C C Jay Kuo, "Videoset: A large-scale compressed video quality dataset based on jnd measurement," *Journal of Visual Communication and Image Representation*, vol. 46, pp. 292–302, 2017.
- [4] Qin Huang, Haiqiang Wang, Sung Chang Lim, Hui Yong Kim, Se Yoon Jeong, and C C Jay Kuo, "Measure and prediction of HEVC perceptually lossy/lossless boundary QP values," *2017 Data Compression Conference (DCC)*, pp. 42–51, 2017.
- [5] Hadi Hadizadeh, Ahmad Reza Heravi, Ivan V Bajic, and Parastoo Karami, "A perceptual distinguishability predictor for JND-noise-contaminated images," *IEEE Trans. Image Process.*, vol. 28, no. 5, pp. 2242–2256, 2019.
- [6] Huanhua Liu, Yun Zhang, Huan Zhang, Chunling Fan, Sam Kwong, C C Jay Kuo, and Xiaoping Fan, "Deep learning-based picture-wise just noticeable distortion prediction model for image compression," *IEEE Trans. Image Process.*, vol. 29, pp. 641–656, 2020.
- [7] Haiqiang Wang, Ioannis Katsavounidis, Qin Huang, Xin Zhou, and C C Jay Kuo, "Prediction of satisfied user ratio for compressed video," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2018, pp. 6747–6751.
- [8] Chunling Fan, Hanhe Lin, Vlad Hosu, Yun Zhang, Qingshan Jiang, Raouf Hamzaoui, and Dietmar Saupe, "SUR-Net: Predicting the satisfied user ratio curve for image compression with deep learning," in *2019 Eleventh International Conference on Quality of Multimedia Experience (QoMEX)*, 2019, pp. 1–6.
- [9] Chunling Fan, Yun Zhang, Huan Zhang, Raouf Hamzaoui, and Qingshan Jiang, "Picture-level just noticeable difference for symmetrically and asymmetrically compressed stereoscopic images: Subjective quality assessment study and datasets," *Journal of Visual Communication and Image Representation*, vol. 62, pp. 140–151, 2019.
- [10] Ming-Jun Chen, Lawrence K. Cormack, and Alan C. Bovik, "No-reference quality assessment of natural stereopairs," *IEEE Trans. Image Process.*, vol. 22, no. 9, pp. 3379–3391, 2013.
- [11] Ming Jun Chen, Che-Chun Su, Do-Kyoung Kwon, Lawrence K. Cormack, and Alan C. Bovik, "Full-reference quality assessment of stereopairs accounting for rivalry," *Signal Processing: Image Communication*, vol. 28, no. 9, pp. 1143–1155, 2013.
- [12] Wenfei Wan, Jinjian Wu, Xuemei Xie, and Guangming Shi, "A novel just noticeable difference model via orientation regularity in DCT domain," *IEEE Access*, vol. 5, pp. 22953–22964, 2017.
- [13] Jinjian Wu, Leida Li, Weisheng Dong, Guangming Shi, Weisi Lin, and C.-C. Jay Kuo, "Enhanced just noticeable difference model for images with pattern complexity," *IEEE Trans. Image Process.*, vol. 26, no. 6, pp. 2682–2693, 2017.
- [14] Zhenyu Wei and King N Ngan, "Spatio-temporal just noticeable distortion profile for grey scale image/video in DCT domain," *IEEE Trans. Circ. Syst. Video Technol.*, vol. 19, no. 3, pp. 337–346, 2009.
- [15] Haiqiang Wang, Ioannis Katsavounidis, Qin Huang, Xin Zhou, and C C Jay Kuo, "Prediction of satisfied user ratio for compressed video," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6747–6751, 2018.
- [16] Sudeng Hu, Lina Jin, Hanli Wang, Yun Zhang, Sam Kwong, and C C Jay Kuo, "Compressed image quality metric based on perceptually weighted distortion," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 5594–5608, 2015.
- [17] David G Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [18] Yun Zhang, Xiangkai Liu, Huanhua Liu, and Chunling Fan, "Depth perceptual quality assessment for symmetrically and asymmetrically distorted stereoscopic 3d videos," *Signal Processing: Image Communication*, vol. 78, pp. 293–305, 2019.
- [19] Alexander J Smola and Bernhard Scholkopf, "A tutorial on support vector regression," *Statistics and Computing*, vol. 14, no. 3, pp. 199–222, 2004.